

Efficient construction of high-resolution TVD conservative schemes for equations with source terms: application to shallow water flows.

J. Burguete and P. Garcia Navarro

Fluid Mechanics. CPS. University of Zaragoza. 50015 Zaragoza. Spain

September 2000

Abstract

High-resolution TVD schemes are widely used for the numerical approximation of conservation laws. The extension to equations with source terms involving spatial derivatives is not immediately apparent. In this paper this extension is defined by limiting the source terms and by including them in the flux limiter functions. On other hand, a manner of build conservative schemes with the non-conservative and with the characteristic forms of the equations is described. In addition, a new treatment of the boundary conditions and a new correction to fix the entropy problem are presented.

Key words

Source-terms, conservative, high-resolution, TVD, flux-limiter, boundaries, entropy

1 Introduction

We are interested in solving as efficiently as possible 1D hyperbolic systems with source terms.

In a general conservative form

$$\frac{\partial \vec{u}(x, t)}{\partial t} + \frac{d\vec{F}(x, \vec{u})}{dx} = \vec{H}(x, \vec{u}) \quad (1.1)$$

where \vec{u} is the vector of conserved variables, \vec{F} the vector of fluxes and \vec{H} that of source terms.

Our interest is led by the numerical modelling of one-dimensional shallow water flows of practical application in Hydraulics such as river flows. In that case

$$\vec{u} = \begin{pmatrix} A \\ Q \end{pmatrix}, \quad \vec{F} = \begin{pmatrix} Q \\ \frac{Q^2}{A} + gI_1 \end{pmatrix}, \quad \vec{H} = \begin{pmatrix} 0 \\ g[I_2 + A(S_0 - S_f)] \end{pmatrix}$$

where Q is the discharge, A is the wetted cross section, g is the acceleration of gravity and S_0 is the bed slope. The rest of the terms account for pressure forces

$$I_1(x, A) = \int_0^{h(x, A)} [h(x, A) - z] \sigma(x, z) dz, \quad I_2(x, A) = \int_0^{h(x, A)} [h(x, A) - z] \frac{\partial \sigma(x, z)}{\partial x} dz$$

(h the water depth and σ the channel width at a position z from the bottom) and for friction forces, with S_f associated to wall friction and represented by the empirical Manning law.

$$S_f = \frac{n^2 Q^2 P^{\frac{4}{3}}}{A^{\frac{10}{3}}}$$

with n the coefficient of wall roughness and P the wetted perimeter.

It is very important to remark that in the conservative form of the equation (1.1) the total derivative $\frac{d\vec{F}}{dx}$ is used since, in this case, the total derivative represents the increments due to the spatial variations in x and in the conserved variable \vec{u} whereas the partial derivative represents the variation due only to the x with \vec{u} constant. The difference between the variations due to both \vec{u} and x from those due to only one of the variables is subtle but significant. Therefore, when the dependence of a general function is $f = f(x, u)$, discrete increments $\frac{\delta f}{\delta x}$ approaches the total derivative $\frac{df}{dx}$, and not the partial derivative $\frac{\partial f}{\partial x}$. The relation between total and partial derivatives is

$$\frac{df(x, u)}{dx} = \frac{\partial f}{\partial x} + \frac{\partial f}{\partial u} \frac{\partial u}{\partial x}$$

From the equations in conservative form (1.1), it is possible to pass to an associated primitive or non-conservative form by noting that

$$\frac{d\vec{F}(x, \vec{u})}{dx} = \frac{\partial \vec{F}(x, \vec{u})}{\partial x} + \frac{\partial \vec{F}(x, \vec{u})}{\partial \vec{u}} \frac{\partial \vec{u}}{\partial x} = \frac{\partial \vec{F}(x, \vec{u})}{\partial x} + \mathbf{J}(x, \vec{u}) \frac{\partial \vec{u}}{\partial x}$$

where $\mathbf{J} = \frac{\partial \vec{F}}{\partial \vec{u}}$ is the Jacobian matrix of the original system. By redefining the source term

$$\vec{H}'(x, \vec{u}) = \vec{H}(x, \vec{u}) - \frac{\partial \vec{F}(x, \vec{u})}{\partial x}$$

so that we get the primitive form:

$$\frac{\partial \vec{u}(x, t)}{\partial t} + \mathbf{J}(x, \vec{u}) \frac{\partial \vec{u}(x, t)}{\partial x} = \vec{H}'(x, \vec{u}) \quad (1.2)$$

Making a carefull derivation, in shallow water equations the following properties are held

$$\begin{aligned} \frac{dh}{dx} &= \frac{\partial h}{\partial x} + \frac{1}{B} \frac{\partial A}{\partial x} \\ \frac{dI_1}{dx} &= \frac{\partial I_1}{\partial x} + \frac{\partial I_1}{\partial A} \frac{\partial A}{\partial x} = I_2 + A \frac{\partial h}{\partial x} + \frac{A}{B} \frac{\partial A}{\partial x} = I_2 + A \frac{dh}{dx} \end{aligned}$$

Then the primitive form of the shallow water equations involves

$$\mathbf{J} = \begin{pmatrix} 0 & 1 \\ c^2 - v^2 & 2v \end{pmatrix}, \quad \vec{H}' = \begin{pmatrix} 0 \\ gA \left[S_0 - S_f - \frac{dh}{dx} + \frac{1}{B} \frac{dA}{dx} \right] \end{pmatrix} \quad (1.3)$$

with B the width at the free surface, $c = \sqrt{g \frac{A}{B}}$ the celerity of infinitesimal surface waves and $v = \frac{Q}{A}$ the fluid velocity. Note that the $\frac{1}{B} \frac{dA}{dx}$ term appears with the carefull derivation and it is not considered in previous works, being very important in the discretizations of the primitive form of the equations.

It is now convenient to develop the characteristic form of the equations given the importance it has for the correct formulation of upwind schemes and boundary conditions. This form is obtained from a diagonalization of the Jacobian in (1.2). Calling \mathbf{P} and \mathbf{P}^{-1} the matrices that make diagonal \mathbf{J} , and $\mathbf{\Lambda}$ the resulting diagonal matrix

$$\mathbf{J} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^{-1}, \quad \mathbf{\Lambda} = \mathbf{P}^{-1} \mathbf{J} \mathbf{P}$$

The matrix $\mathbf{\Lambda}$ is formed by the eigenvalues of \mathbf{J} , and \mathbf{P} is constructed with its eigenvectors. Let \vec{w} be the set of variables (characteristic variables) that verify

$$d\vec{u} = \mathbf{P} d\vec{w}, \quad d\vec{w} = \mathbf{P}^{-1} d\vec{u}$$

Then,

$$\frac{\partial \vec{w}(x, t)}{\partial t} + \mathbf{\Lambda}(x, \vec{w}) \frac{\partial \vec{w}(x, t)}{\partial x} = \mathbf{P}^{-1}(x, \vec{w}) \vec{H}'(x, \vec{w}) \quad (1.4)$$

Note that the characteristic variables as defined are not functions since, in general,

$$\frac{\partial w_i}{\partial u_j \partial u_k} \neq \frac{\partial w_i}{\partial u_k \partial u_j}$$

In the shallow water equations, the above matrices are

$$\mathbf{P} = \begin{pmatrix} 1 & 1 \\ v+c & v-c \end{pmatrix}, \quad \mathbf{P}^{-1} = \frac{1}{2c} \begin{pmatrix} c-v & 1 \\ c+v & -1 \end{pmatrix}, \quad \mathbf{\Lambda} = \begin{pmatrix} v+c & 0 \\ 0 & v-c \end{pmatrix}$$

2 Conservative schemes

The conservation law (1.1) contains an important physical meaning. By spatial integration

$$\int_0^L \left(\frac{\partial \vec{u}}{\partial t} + \frac{d\vec{F}}{dx} \right) dx = \int_0^L \vec{H} dx \Rightarrow \int_0^L \frac{\partial \vec{u}}{\partial t} dx = \vec{F}_0 - \vec{F}_L + \int_0^L \vec{H} dx \quad (2.1)$$

It is expressed that the time variation in the conserved variable in a given volume is equal to the difference between the incoming and the outgoing fluxes plus the contribution of the source term. When discretizing a conservation law of this kind, bad numerical approximations can lead to bad behaviour in the solution and unacceptable error. Schemes properly approximating the conservation equation (2.1) are called conservative schemes. Some ways of defining them are presented next.

2.1 Conservative schemes with numerical flux

The most common definition of a conservative scheme follows the structure

$$\frac{\Delta \vec{u}_i^n}{\Delta t} = \vec{H}_i^* - \frac{1}{\delta x} \left(\vec{F}_{i+\frac{1}{2}}^* - \vec{F}_{i-\frac{1}{2}}^* \right) \quad (2.2)$$

where \vec{H}^* and \vec{F}^* are the numerical source and the numerical flux respectively and represent a suitable approximation to the true source and flux terms in the equation so that the scheme gets the properties required. Δ will be used for time increments $\Delta f_i^n = f_i^{n+1} - f_i^n$, and δ represents spatial increment $\delta f_{i+\frac{1}{2}}^n = f_{i+1}^n - f_i^n$. Schemes so defined will be conservative since they produce a good approximation of (2.1) canceling the contributions of the flux at the intermedial interfaces, being the variation of the conserved variable due only to the source terms and to the flux at the

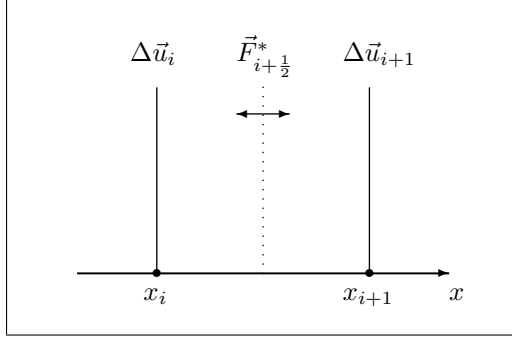


Figure 2.1: Adding the variable increments $\Delta \vec{u}_i$ the numerical flux contributions in the intermedial interfaces $\vec{F}_{i+\frac{1}{2}}^*$ are cancelled producing a conservative scheme.

boundaries

$$\sum_{i=1}^{N-1} \frac{\Delta \vec{u}_i^n}{\Delta t} \delta x = \frac{\Delta}{\Delta t} \sum_{i=1}^{N-1} \vec{u}_i^n \delta x \approx \frac{\partial}{\partial t} \int_{x_{\frac{1}{2}}}^{x_{N-\frac{1}{2}}} \vec{u} dx$$

$$\sum_{i=1}^{N-1} [\vec{H}_i^* \delta x - (\vec{F}_{i+\frac{1}{2}}^* - \vec{F}_{i-\frac{1}{2}}^*)] = \sum_{i=1}^{N-1} [\vec{H}_i^* \delta x] + \vec{F}_{\frac{1}{2}}^* - \vec{F}_{N-\frac{1}{2}}^* \approx \vec{F}_{\frac{1}{2}} - \vec{F}_{N-\frac{1}{2}} + \int_{x_{\frac{1}{2}}}^{x_{N-\frac{1}{2}}} \vec{H} dx$$

2.2 Conservative schemes with wave decomposition

A total numerical flux \vec{F}_i^T can also be defined at the grid nodes. The difference in this flux across a grid cell can be decomposed into incoming $(\delta \vec{F}_{i+\frac{1}{2}}^R)$ and outgoing $(\delta \vec{F}_{i+\frac{1}{2}}^L)$ parts. Schemes so built follow

$$\delta \vec{F}_{i+\frac{1}{2}}^T = \vec{F}_{i+1}^T - \vec{F}_i^T = \delta \vec{F}_{i+\frac{1}{2}}^R + \delta \vec{F}_{i+\frac{1}{2}}^L$$

$$\frac{\Delta \vec{u}_i^n}{\Delta t} = \vec{H}_i^* - \frac{1}{\delta x} (\delta \vec{F}_{i+\frac{1}{2}}^R + \delta \vec{F}_{i-\frac{1}{2}}^L) \quad (2.3)$$

This also leads to conservative schemes since this form can be shown to be equivalent to (2.2). It

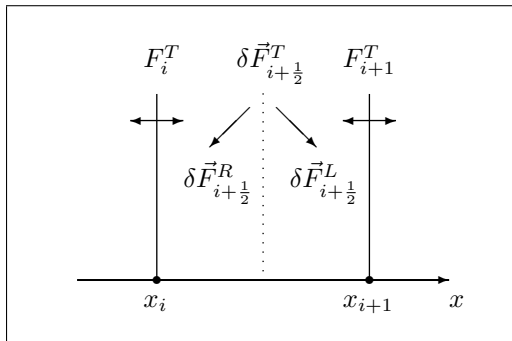


Figure 2.2: Adding the total flux contributions $\delta \vec{F}_{i+\frac{1}{2}}^R$ and $\delta \vec{F}_{i+\frac{1}{2}}^L$ the flux contributions in the internal grid cells \vec{F}_i^T are cancelled producing a conservative scheme.

is possible to define a function $\vec{\Phi}$ so that

$$\delta \vec{F}_{i+\frac{1}{2}}^L = \frac{1}{2} \delta \vec{F}_{i+\frac{1}{2}}^T + \vec{\Phi}_{i+\frac{1}{2}}, \quad \delta \vec{F}_{i+\frac{1}{2}}^R = \frac{1}{2} \delta \vec{F}_{i+\frac{1}{2}}^T - \vec{\Phi}_{i+\frac{1}{2}}$$

So that from (2.3)

$$\frac{\Delta \vec{u}_i^n}{\Delta t} = \vec{H}_i^* - \frac{1}{\delta x} \left[\frac{1}{2} (\vec{F}_{i+1}^T - \vec{F}_i^T) - \vec{\Phi}_{i+\frac{1}{2}} + \frac{1}{2} (\vec{F}_i^T - \vec{F}_{i-1}^T) + \vec{\Phi}_{i-\frac{1}{2}} \right]$$

and the following interface numerical flux can be defined:

$$\vec{F}_{i+\frac{1}{2}}^* = \frac{1}{2} (\vec{F}_{i+1}^T + \vec{F}_i^T) - \vec{\Phi}_{i+\frac{1}{2}} = \vec{F}_i^T + \delta \vec{F}_{i+\frac{1}{2}}^R = \vec{F}_{i+1}^T - \delta \vec{F}_{i+\frac{1}{2}}^L \quad (2.4)$$

In addition, it can be interesting to consider a non-centered contribution of the source terms

$$\begin{aligned} \vec{H}_{i+\frac{1}{2}}^T &= \vec{H}_{i+\frac{1}{2}}^R + \vec{H}_{i+\frac{1}{2}}^L \\ \frac{\Delta \vec{u}_i^n}{\Delta t} &= \left(\vec{H} - \frac{\delta \vec{F}}{\delta x} \right)_{i-\frac{1}{2}}^L + \left(\vec{H} - \frac{\delta \vec{F}}{\delta x} \right)_{i+\frac{1}{2}}^R \end{aligned} \quad (2.5)$$

Some proofs of the advantage of the source terms decentralisation are shown in other works [3,4].

In order to illustrate the mean of this wave decomposition, two examples are here mentioned.

First order upwind discretizations admits a decomposition like

$$\vec{F}_i^T = \vec{F}_i^n, \quad \delta \vec{F}_{i+\frac{1}{2}}^L = \left(\delta \vec{F}^+ \right)_{i+\frac{1}{2}}^n, \quad \delta \vec{F}_{i+\frac{1}{2}}^R = \left(\delta \vec{F}^- \right)_{i+\frac{1}{2}}^n$$

where $\delta \vec{F}^+$ and $\delta \vec{F}^-$ associated to positive and negative propagation velocities respectively, whereas

first order in time centered discretization admits a decomposition like

$$\vec{F}_i^T = \vec{F}_i^n, \quad \delta \vec{F}_{i+\frac{1}{2}}^L = \delta \vec{F}_{i+\frac{1}{2}}^R = \frac{1}{2} \delta \vec{F}_{i+\frac{1}{2}}^n$$

2.3 Conservative schemes in primitive form

Conservative schemes can also be derived from the primitive form of the equations (1.2). The advantage is that the latter form tends to be simpler to deal with than the conservative form.

We need to establish the conditions under which schemes derived this way are equivalent to the conservative schemes derived from the conservative equations. First of all, the following equality must hold at the discrete level

$$\vec{G}_{i+1/2} \equiv \left(\vec{H} - \frac{\delta \vec{F}}{\delta x} \right)_{i+\frac{1}{2}} = \left(\vec{H}' - \mathbf{J} \frac{\delta \vec{u}}{\delta x} \right)_{i+\frac{1}{2}} \quad (2.6)$$

Note that this equality requires a non-pointwise treatment of source terms and is equivalent, with any detail due to the source terms, to the Roe's average [8,11]. To make a primitive scheme with this average ensures the scheme to be conservative. Once this has been achieved, two equivalent forms to build conservative schemes with decentered source terms are possible. Defining the generalised source term \vec{G} like

$$\vec{G}_{i+1/2} \equiv \left(\vec{H} - \frac{\delta \vec{F}}{\delta x} \right)_{i+\frac{1}{2}} \quad (2.7)$$

the standard conservative form is achieved, whereas defining this term like

$$\vec{G}_{i+1/2} \equiv \left(\vec{H}' - \mathbf{J} \frac{\delta \vec{u}}{\delta x} \right)_{i+\frac{1}{2}} \quad (2.8)$$

with the restriction (2.6), the conservative in primitive form version of the scheme is gotten. Then, a wave decomposition equivalent to (2.5) can be performed

$$\frac{\Delta \vec{u}_i^n}{\Delta t} = \vec{G}_{i-\frac{1}{2}}^L + \vec{G}_{i+\frac{1}{2}}^R \quad (2.9)$$

Condition (2.6) imposes in the shallow water equations the following equality

$$\begin{aligned} \left(\begin{array}{c} 0 \\ g[I_2 + A(S_0 - S_f)] \end{array} \right)_{i+\frac{1}{2}} - \frac{\delta}{\delta x} \left(\begin{array}{c} Q \\ \frac{Q^2}{A} + gI_1 \end{array} \right)_{i+\frac{1}{2}} &= \left(\begin{array}{c} 0 \\ gA(S_0 - S_f + \frac{1}{B} \frac{dA}{dx} - \frac{dh}{dx}) \end{array} \right)_{i+\frac{1}{2}} \\ &- \left(\begin{array}{cc} 0 & 1 \\ c^2 - v^2 & 2v \end{array} \right)_{i+\frac{1}{2}} \frac{\delta}{\delta x} \left(\begin{array}{c} A \\ Q \end{array} \right)_{i+\frac{1}{2}} \end{aligned}$$

which, considering the following discretizations of second order of approximation

$$\left(\frac{dh}{dx} \right)_{i+\frac{1}{2}} \approx \left(\frac{\delta h}{\delta x} \right)_{i+\frac{1}{2}}, \quad \left(\frac{1}{B} \frac{dA}{dx} \right)_{i+\frac{1}{2}} \approx \left(\frac{1}{B} \frac{\delta A}{\delta x} \right)_{i+\frac{1}{2}}, \quad (I_2)_{i+\frac{1}{2}} \approx \left(\frac{\delta I_1}{\delta x} - A \frac{\delta h}{\delta x} \right)_{i+\frac{1}{2}}$$

it reduces to

$$-\frac{\delta}{\delta x} \left(\frac{Q^2}{A} \right)_{i+\frac{1}{2}} = \left[g \frac{A}{B} \frac{\delta A}{\delta x} - (c^2 - v^2) \frac{\delta A}{\delta x} - 2v \frac{\delta Q}{\delta x} \right]_{i+\frac{1}{2}}$$

This condition can be shown to hold with the following average values

$$c_{i+\frac{1}{2}} = \sqrt{g \frac{A_{i+\frac{1}{2}}}{B_{i+\frac{1}{2}}}}, \quad v_{i+\frac{1}{2}} = \frac{Q_{i+1}/\sqrt{A_{i+1}} + Q_i/\sqrt{A_i}}{\sqrt{A_{i+1}} + \sqrt{A_i}} \quad (2.10)$$

The choice of the discrete averages $A_{i+\frac{1}{2}}$, $B_{i+\frac{1}{2}}$, $(S_0)_{i+\frac{1}{2}}$ and $(S_f)_{i+\frac{1}{2}}$ is open but, in our work, this has not proved significant. The simplest arithmetic average has been applied.

2.4 Conservative schemes in characteristic form

It is also possible to derive conservative schemes based in the characteristic form of the equations. This will be the basis for the wave decomposition of the upwind schemes. From (1.4) it is possible to rewrite

$$\frac{\partial \vec{w}}{\partial t} = \mathbf{P}^{-1} \left(\vec{H}' - \mathbf{J} \frac{\partial \vec{u}}{\partial x} \right)$$

Next, a discrete wave decomposition into left and right moving contributions can be done

$$\left(\mathbf{P}^{-1} \vec{G} \right)_{i+\frac{1}{2}} = \left(\boldsymbol{\Omega}^L \mathbf{P}^{-1} \vec{G} \right)_{i+\frac{1}{2}} + \left(\boldsymbol{\Omega}^R \mathbf{P}^{-1} \vec{G} \right)_{i+\frac{1}{2}}$$

being $\boldsymbol{\Omega}^L$ and $\boldsymbol{\Omega}^R$ diagonal matrices to be defined in every particular numerical scheme. In order to ensure the conservative character of the scheme, they have to obey

$$(\boldsymbol{\Omega}^L + \boldsymbol{\Omega}^R)_{i+\frac{1}{2}} = \mathbf{I} \quad (2.11)$$

Returning to the physical variable by multiplying by \mathbf{P} , the final form for the discretization is

$$\frac{\Delta \vec{u}_i^n}{\Delta t} = \left(\mathbf{P} \boldsymbol{\Omega}^L \mathbf{P}^{-1} \vec{G} \right)_{i-\frac{1}{2}} + \left(\mathbf{P} \boldsymbol{\Omega}^R \mathbf{P}^{-1} \vec{G} \right)_{i+\frac{1}{2}} \quad (2.12)$$

Note that this discretization requires again a decentralised formulation of source terms, being equally possible a conservative or a primitive definition for \vec{G} . If a pointwise treatment of the source terms is desired, it can be made

$$\frac{\Delta \vec{u}_i^n}{\Delta t} = \vec{H}_i^n - \left(\mathbf{P} \boldsymbol{\Omega}^L \mathbf{P}^{-1} \frac{\delta \vec{F}}{\delta x} \right)_{i-\frac{1}{2}} - \left(\mathbf{P} \boldsymbol{\Omega}^R \mathbf{P}^{-1} \frac{\delta \vec{F}}{\delta x} \right)_{i+\frac{1}{2}} \quad (2.13)$$

3 First order upwind scheme

Upwind schemes are based on the idea of approximating the spatial derivatives by non-centered differences biased in the sense of propagation of information in the physical problem. In order to construct a first order scheme, suitable for left and right moving propagation velocities, the following can be written

$$\Delta \vec{u}_i^n = \Delta t \left[\vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] \quad (3.1)$$

where $\delta \vec{F}^+$ is associated to negative velocities and $\delta \vec{F}^-$ to positive velocities. A linear analysis of the homogeneous equations shows that the stability condition is $CFL \leq 1$ and it will be dissipative

provided that $CFL < 1$, with $CFL = \max |a_k| \frac{\Delta t}{\delta x}$ the CFL number and a_k being the eigenvalues of the Jacobian.

When the source terms are dominant in a problem, it may be necessary to introduce a semi-implicit treatment for them in order to stabilise the scheme. One way to proceed is

$$\vec{H}_i^* = \theta \vec{H}_i^{n+1} + (1 - \theta) \vec{H}_i^n \approx \vec{H}_i^n + \theta \left(\frac{\partial \vec{H}}{\partial \vec{u}} \frac{\partial \vec{u}}{\partial t} \right)_i^n \Delta t = \vec{H}_i^n + \theta \mathbf{K}_i^n \Delta \vec{u}_i^n$$

where $\mathbf{K} = \frac{\partial \vec{H}}{\partial \vec{u}}$ is the Jacobian of the source term. Putting this into (3.1), the scheme becomes:

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \Delta t \left[\vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right]$$

With the addition of the source terms in semi-implicit form, (2.12) can be written

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \Delta t \left[\left(\mathbf{P} \mathbf{\Omega}^L \mathbf{P}^{-1} \vec{G} \right)_{i-\frac{1}{2}}^n + \left(\mathbf{P} \mathbf{\Omega}^R \mathbf{P}^{-1} \vec{G} \right)_{i+\frac{1}{2}}^n \right]$$

It will maintain the conservative properties with the restriction (2.11). Furthermore, the following wave decomposition is assumed in order to select the appropriate influence region in every case.

$$\mathbf{\Omega}^L = \mathbf{\Omega}^+ = \frac{1}{2} [\mathbf{I} + \text{sign}(\mathbf{\Lambda})], \quad \mathbf{\Omega}^R = \mathbf{\Omega}^- = \frac{1}{2} [\mathbf{I} - \text{sign}(\mathbf{\Lambda})] \quad (3.2)$$

so that:

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \Delta t \left[\left(\mathbf{P} \mathbf{\Omega}^+ \mathbf{P}^{-1} \vec{G} \right)_{i-\frac{1}{2}}^n + \left(\mathbf{P} \mathbf{\Omega}^- \mathbf{P}^{-1} \vec{G} \right)_{i+\frac{1}{2}}^n \right]$$

With the notation

$$\vec{G}^\pm = \mathbf{P} \mathbf{\Omega}^\pm \mathbf{P}^{-1} \vec{G} \quad (3.3)$$

the scheme translates back to a simpler form

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \Delta t \left[\left(\vec{G}^+ \right)_{i-\frac{1}{2}}^n + \left(\vec{G}^- \right)_{i+\frac{1}{2}}^n \right] \quad (3.4)$$

This numerical scheme as defined in (3.4) does not produce good results in the presence of transcritical flow. This corresponds to a change of sign from negative to positive in the advection velocities. The numerical scheme is not able to interpret the transition as smooth and gives rise to non-physical shocks (entropy problem). This situation is associated to a local lack of numerical dissipation that can be corrected as outlined in Appendix A.

It is worth stressing here that this first order upwind scheme produces second order accuracy in space for steady cases. Assuming for simplicity a positive propagation velocities, eliminating the time dependence from the scheme (3.4) we get

$$\left(\vec{H} - \frac{\delta \vec{F}}{\delta x} \right)_{i-\frac{1}{2}}^n = 0$$

Which means that

$$\vec{F}_i^n = \vec{F}_{i-1}^n + \delta x \vec{H}_{i-\frac{1}{2}}^n$$

This is the mean point integration rule, an approximation of second order. This gain in accuracy is the main advantage of the upwinding of the source terms. Actually, taking a pointwise approach

$$\vec{F}_i^n = \vec{F}_{i-1}^n + \delta x \vec{H}_i^n$$

which is Euler integration rule, a first order approximation.

The conservative character of this scheme is proved by the existence of a wave decomposition like

$$\vec{F}_i^T = \vec{F}_i^n, \quad \delta \vec{F}_{i+\frac{1}{2}}^L = \left(\delta \vec{F}^+ \right)_{i+\frac{1}{2}}^n, \quad \delta \vec{F}_{i+\frac{1}{2}}^R = \left(\delta \vec{F}^- \right)_{i+\frac{1}{2}}^n$$

and a cell numerical flux:

$$\vec{F}_{i+\frac{1}{2}}^* = \frac{1}{2} \left\{ \vec{F}_i^n + \vec{F}_{i+1}^n - \left[\mathbf{P} \text{sign}(\mathbf{\Lambda}) \mathbf{P}^{-1} \delta \vec{F} \right]_{i+\frac{1}{2}}^n \right\}$$

4 Spatially second order TVD scheme

4.1 Second order in space TVD scheme for the scalar case

As a preliminary step before the statement of the form of this functions, the TVD conditions will be outlined for the scalar case. Assume for simplicity a homogeneous scalar conservation law

$$\frac{\partial u(x, t)}{\partial t} + \frac{\partial F(u)}{\partial x} = 0 \Rightarrow \quad \frac{\partial u(x, t)}{\partial t} + a(u) \frac{\partial u(x, t)}{\partial x} = 0$$

with $a = \frac{\partial F}{\partial u}$. The Total Variation of the discrete solution to this equation is defined:

$$TV^n = \sum_i |u_{i+1}^n - u_i^n|$$

Then a scheme will be TVD if it obeys the following property that prevents oscillatory solutions

$$TV^{n+1} = \sum_i |u_{i+1}^{n+1} - u_i^{n+1}| \leq \sum_i |u_{i+1}^n - u_i^n| = TV^n \quad (4.1)$$

A general explicit scheme that can be written in the form

$$\Delta u_i^n = -\lambda(\delta F_{i-\frac{1}{2}}^+ + \delta F_{i+\frac{1}{2}}^-)$$

with $\lambda = \frac{\Delta t}{\Delta x}$, and it admits the definition of the following coefficients

$$\delta F_{i+\frac{1}{2}}^+ = C_{i+\frac{1}{2}}^+ \delta u_{i+\frac{1}{2}}^n, \quad \delta F_{i+\frac{1}{2}}^- = C_{i+\frac{1}{2}}^- \delta u_{i+\frac{1}{2}}^n$$

The scheme will be TVD if, using (4.1) (see [5] for more detail),

$$C_{i+\frac{1}{2}}^- \leq 0, \quad C_{i+\frac{1}{2}}^+ \geq 0, \quad \lambda(C_{i+\frac{1}{2}}^+ - C_{i+\frac{1}{2}}^-) \leq 1 \quad (4.2)$$

This conditions are automatically fulfilled by the first order upwind scheme with the *CFL* condition.

The scheme for the scalar case involves scalar flux limiter functions

$$\begin{aligned} \Delta \tilde{u}_i^n = -\lambda \Big\{ & (\delta F^+)^n_{i-\frac{1}{2}} + (\delta F^-)^n_{i+\frac{1}{2}} + \frac{1}{2} \left[(\Psi^+ \delta F^+)^n_{i-\frac{1}{2}} - (\Psi^+ \delta F^+)^n_{i-\frac{3}{2}} \right. \\ & \left. + (\Psi^- \delta F^-)^n_{i+\frac{1}{2}} - (\Psi^- \delta F^-)^n_{i+\frac{3}{2}} \right] \Big\} \end{aligned}$$

Then the flux limiter functions are defined to combine the second order spatial centered and upwind schemes, for preserving the second order, and according to the (4.2) properties, for avoiding the numerical oscillations. In order to produce a second order scheme, the dependence of flux limiter functions is defined like:

$$\Psi_{i+\frac{1}{2}}^+ = \Psi(r_{i+\frac{1}{2}}^+) = \Psi\left(\frac{\delta F_{i+\frac{3}{2}}^+}{\delta F_{i+\frac{1}{2}}^+}\right), \quad \Psi_{i+\frac{1}{2}}^- = \Psi(r_{i+\frac{1}{2}}^-) = \Psi\left(\frac{\delta F_{i-\frac{1}{2}}^-}{\delta F_{i+\frac{1}{2}}^-}\right),$$

With this dependence, the following properties are achieved:

$\Psi(r) = 1, \forall r \Rightarrow$ Upwind scheme second order accurate in space (unstable)

$\Psi(r) = r, \forall r \Rightarrow$ Central Lax's scheme second order accurate in space (unstable)

Infinite spatially second order schemes can be created for intermediate values between $\Psi(r) = 1$ and $\Psi(r) = r$. The zone of second order in space is shaded in Fig. 4.1. Furthermore the TVD properties

are required to create a scheme without numerical oscillations. Applying the TVD conditions (4.2) the flux limiter will be a positive function so that

$$\Psi(r) = 0, \forall r < 0; \quad \Psi(r) \leq 2r, \forall r > 0$$

and it produces the following stability condition:

$$CFL \leq \frac{1}{1 + \frac{1}{2} \max(\Psi)} \quad (4.3)$$

For working with reasonable time steps it is usual to establish the restriction $\Psi(r) \leq 2$ in order to be capable to work up to $CFL \leq \frac{1}{2}$. The intersection between the second order region and the TVD region for the flux limiter functions in the second order in space TVD scheme is represented in the Fig. 4.2. Many particular flux limiter functions are defined in other works. We use the most usual:

Superbee: $\Psi(r) = \max[0, \min(1, 2r), \min(2, r)]$

Van-Leer: $\Psi(r) = \frac{r+|r|}{1+|r|}$

Van-Albada: $\Psi(r) = \begin{cases} \frac{r+r^2}{1+r^2}, & \forall r > 0 \\ 0, & \forall r \leq 0 \end{cases}$

Minmod: $\Psi(r) = \max[0, \min(1, r)]$

4.2 Extension to systems with source terms

The unstable Lax's scheme consists of a second order central approximation to the spatial derivative and a first order approximation to the time derivative leading to the scheme:

$$\Delta \vec{u}_i^n = \Delta t \left[\vec{H}_i^n - \frac{1}{2\delta x} (\vec{F}_{i+1}^n - \vec{F}_{i-1}^n) \right]$$

or, as a combination of the first order upwind (3.1) plus second order terms:

$$\begin{aligned} \Delta \vec{u}_i^n = \Delta t & \left[\vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] + \frac{\Delta t}{2} \left[\left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n + \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right. \\ & \left. - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i+\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i-\frac{1}{2}}^n \right] \end{aligned} \quad (4.4)$$

It is equally possible to arrive to an upwind scheme with similar properties by means of a second order upwind approximation to the space derivative [5]. By doing so, considering the possibility of both positive and negative advection velocities,

$$\Delta \vec{u}_i^n = \Delta t \left\{ \vec{H}_i^n - \frac{3}{2} \left[\left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n + \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] + \frac{1}{2} \left[\left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{3}{2}}^n + \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{3}{2}}^n \right] \right\}$$

or, in terms of the first order upwind

$$\Delta \vec{u}_i^n = \Delta t \left[\vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] + \frac{\Delta t}{2} \left[\left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{3}{2}}^n + \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{3}{2}}^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] \quad (4.5)$$

These two schemes (4.4) and (4.5), although unstable, are the basis for the construction of the second order in space TVD schemes. By means of an adequate limitation of the spatial second order terms, second order accuracy can be preserved whilst oscillations are avoided. Multiplying the flux differences by the limiting functions [5] and including a semi-implicit treatment of the source terms,

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \Delta t \left\{ \vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n + \frac{1}{2} \left[\left(\Psi^+ \frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{3}{2}}^n + \left(\Psi^- \frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{3}{2}}^n - \left(\Psi^+ \frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\Psi^- \frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] \right\} \quad (4.6)$$

where Ψ^+ and Ψ^- are in general matrices.

Once the conditions on the flux limiter functions have been established so that the desired properties are achieved on the solution to a scalar equation, a generalization to systems with source terms is desired. One of the simplest ways to define the flux limiting matrices Ψ^+ y Ψ^- for a flux $\vec{F} = (F^1, \dots, F^k)$ is

$$\Psi_{i+\frac{1}{2}}^\pm = \begin{pmatrix} \Psi \left(\frac{(\delta F^1)^\pm_{i+\frac{1}{2} \pm 1}}{(\delta F^1)^\pm_{i+\frac{1}{2}}} \right) & & \\ & \ddots & \\ & & \Psi \left(\frac{(\delta F^k)^\pm_{i+\frac{1}{2} \pm 1}}{(\delta F^k)^\pm_{i+\frac{1}{2}}} \right) \end{pmatrix} \quad (4.7)$$

Another form to achieve this is shown in [1].

This TVD scheme second order in space is conservative since it admits the following wave decomposition:

$$\begin{aligned} \vec{F}_i^T &= \vec{F}_i^n \\ \delta \vec{F}_{i+\frac{1}{2}}^L &= (\delta \vec{F}^+)_{i+\frac{1}{2}}^n - \frac{1}{2} (\Psi^+ \delta \vec{F}^+)_{i-\frac{1}{2}}^n + \frac{1}{2} (\Psi^- \delta \vec{F}^-)_{i+\frac{3}{2}}^n \\ \delta \vec{F}_{i+\frac{1}{2}}^R &= (\delta \vec{F}^-)_{i+\frac{1}{2}}^n + \frac{1}{2} (\Psi^+ \delta \vec{F}^+)_{i-\frac{1}{2}}^n - \frac{1}{2} (\Psi^- \delta \vec{F}^-)_{i+\frac{3}{2}}^n \end{aligned}$$

and the following numerical flux:

$$\vec{F}_{i+\frac{1}{2}}^* = \frac{1}{2} \left\{ \vec{F}_i^n + \vec{F}_{i+1}^n - \left[\mathbf{P} \text{sign}(\mathbf{\Lambda}) \mathbf{P}^{-1} \delta \vec{F} \right]_{i+\frac{1}{2}}^n + \left(\mathbf{\Psi}^+ \delta \vec{F}^+ \right)_{i-\frac{1}{2}}^n - \left(\mathbf{\Psi}^- \delta \vec{F}^- \right)_{i+\frac{3}{2}}^n \right\}$$

Our interest is focussed on the correct treatment of problems involving source terms. We will do that working with the already defined generalised source term $\vec{G}_{i+\frac{1}{2}}$. With this notation, the conservative TVD second order in space scheme can be written

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \Delta t \left\{ \left(\vec{G}^+ \right)_{i-\frac{1}{2}}^n + \left(\vec{G}^- \right)_{i+\frac{1}{2}}^n + \frac{1}{2} \left[\left(\mathbf{\Psi}^+ \vec{G}^+ \right)_{i-\frac{1}{2}}^n + \left(\mathbf{\Psi}^- \vec{G}^- \right)_{i+\frac{1}{2}}^n - \left(\mathbf{\Psi}^+ \vec{G}^+ \right)_{i-\frac{3}{2}}^n - \left(\mathbf{\Psi}^- \vec{G}^- \right)_{i+\frac{3}{2}}^n \right] \right\} \quad (4.8)$$

Being the unstable Lax's scheme with non-centered source term

$$(1 - \theta \mathbf{K}_i^n \Delta t) \Delta \vec{u}_i^n = \frac{\Delta t}{2} \left[\left(\vec{G}^+ \right)_{i-\frac{1}{2}}^n + \left(\vec{G}^- \right)_{i-\frac{1}{2}}^n + \left(\vec{G}^+ \right)_{i+\frac{1}{2}}^n + \left(\vec{G}^- \right)_{i+\frac{1}{2}}^n \right]$$

In order to ensure this limit when $\Psi(r) = r$ the limiting diagonal matrices $\mathbf{\Psi}^+$ and $\mathbf{\Psi}^-$, for

$\vec{G} = (G^1, \dots, G^k)$, must be

$$\mathbf{\Psi}_{i+\frac{1}{2}}^\pm = \begin{pmatrix} \Psi \left(\frac{(G^1)_{i+\frac{1}{2}}^\pm}{(G^1)_{i+\frac{1}{2}}^\pm} \right) & & \\ & \ddots & \\ & & \Psi \left(\frac{(G^k)_{i+\frac{1}{2}}^\pm}{(G^k)_{i+\frac{1}{2}}^\pm} \right) \end{pmatrix} \quad (4.9)$$

It is worth stressing that a rigorous upwind treatment of the source terms in TVD schemes imposes not only the limitation of the source terms but also their involvement in the definition of the limitation function itself. However it is possible to neglect the contributions of the source terms in the limiter functions using the limiter matrix (4.7) but this, although the limitation of the flux is ensured and it makes a scheme without numerical oscillations, do not ensures the second order of the scheme and a loss of accuracy is produced. Furthermore this is not necessary because the vector components G^k must be equally calculated and it makes simpler to compute the limiter matrix (4.9).

5 Second order in space and time TVD scheme

5.1 Second order in space and time schemes

To develop the schemes to second order in time implies the approximation of the time derivatives to second order, which, using a system of equations like (1.1)

$$\Delta \vec{u}_i^n = \left(\frac{\partial \vec{u}}{\partial t} \right)_i^n \Delta t + \frac{1}{2} \left(\frac{\partial^2 \vec{u}}{\partial t^2} \right)_i^n \Delta t^2 + O(\Delta t^3) = \left(\vec{H} - \frac{\partial \vec{F}}{\partial x} \right)_i^n \Delta t + \frac{1}{2} \frac{\partial}{\partial t} \left(\vec{H} - \frac{\partial \vec{F}}{\partial x} \right)_i^n \Delta t^2 + O(\Delta t^3)$$

For the time derivative the following has to be taken into account

$$\frac{\partial}{\partial t} \left(\vec{H} - \frac{\partial \vec{F}}{\partial x} \right) = \frac{\partial \vec{H}}{\partial \vec{u}} \frac{\partial \vec{u}}{\partial t} - \frac{\partial}{\partial x} \left(\frac{\partial \vec{F}}{\partial \vec{u}} \frac{\partial \vec{u}}{\partial t} \right) = \mathbf{K} \frac{\partial \vec{u}}{\partial t} - \frac{\partial}{\partial x} \left[\mathbf{J} \left(\vec{H} - \frac{\partial \vec{F}}{\partial x} \right) \right]$$

where the Jacobians of the flux and source term have been used. Remark that in this second order term it appears $\mathbf{K} \frac{\partial \vec{u}}{\partial t}$ and $\frac{\partial(\mathbf{J}\vec{H})}{\partial x}$ terms which are due to the presence of the source term in the equation. This terms are obviated in previous works. Making a central approximation of the spatial derivatives

$$\frac{\partial}{\partial t} \left(\vec{H} - \frac{\partial \vec{F}}{\partial x} \right)_i^n = \mathbf{K}_i^n \frac{\Delta \vec{u}_i^n}{\Delta t} - \frac{1}{\delta x} \left[\left(\mathbf{J}\vec{G} \right)_{i+\frac{1}{2}}^n - \left(\mathbf{J}\vec{G} \right)_{i-\frac{1}{2}}^n \right] + O(\delta x, \Delta t)$$

where, for simplicity and stability, a semi-implicit treatment of the source term is made. The Lax-Wendroff scheme with central source terms can then be expressed

$$\left(1 - \mathbf{K}_i^n \frac{\Delta t}{2} \right) \Delta \vec{u}_i^n = \Delta t \left(\vec{H}_i^n - \frac{\vec{F}_{i+1}^n - \vec{F}_i^n}{2\delta x} \right) - \frac{\Delta t^2}{2\delta x} \left[\left(\mathbf{J}\vec{G} \right)_{i+\frac{1}{2}}^n - \left(\mathbf{J}\vec{G} \right)_{i-\frac{1}{2}}^n \right] \quad (5.1)$$

and represents a central second order in space and time approximation. This scheme is stable provided that $CFL \leq 1$ and dissipative if $CFL < 1$. Written as a sum of the first order upwind plus second order correction terms:

$$\begin{aligned} \left(1 - \mathbf{K}_i^n \frac{\Delta t}{2} \right) \Delta \vec{u}_i^n = & \left[\vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] + \frac{\Delta t}{2\delta x} \left[\left(\delta \vec{F}^- - \Delta t \mathbf{J}^- \vec{G}^- \right)_{i+\frac{1}{2}}^n \right. \\ & \left. - \left(\delta \vec{F}^- - \Delta t \mathbf{J}^- \vec{G}^- \right)_{i-\frac{1}{2}}^n - \left(\delta \vec{F}^+ + \Delta t \mathbf{J}^+ \vec{G}^+ \right)_{i+\frac{1}{2}}^n + \left(\delta \vec{F}^+ + \Delta t \mathbf{J}^+ \vec{G}^+ \right)_{i-\frac{1}{2}}^n \right] \end{aligned}$$

To arrive to an upwind second order in space and time approximation, the starting point is the same but the space derivatives are approximated in a non-centered form. Ensuring that it is properly defined for both positive and negative senses of propagation, the following scheme is obtained

$$\left(1 - \mathbf{K}_i^n \frac{\Delta t}{2} \right) \Delta \vec{u}_i^n = \Delta t \left\{ \vec{H}_i^n - \frac{3}{2} \left[\left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n + \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] + \frac{1}{2} \left[\left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{3}{2}}^n \right. \right.$$

$$+ \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{3}{2}}^n \Big] - \frac{\Delta t^2}{2\delta x} \left[\left(\mathbf{J}^- \vec{G}^- \right)_{i+\frac{3}{2}}^n - \left(\mathbf{J}^- \vec{G}^- \right)_{i+\frac{1}{2}}^n + \left(\mathbf{J}^+ \vec{G}^+ \right)_{i-\frac{1}{2}}^n - \left(\mathbf{J}^+ \vec{G}^+ \right)_{i-\frac{3}{2}}^n \right] \Big\} \quad (5.2)$$

The stability condition for this scheme is $CFL \leq 2$ and it is dissipative if $CFL < 2$.

A TVD scheme second order accurate in space and time can be built from (5.2) as a sum of the first order upwind plus limited second order correction terms, as

$$\begin{aligned} \left(1 - \mathbf{K}_i^n \frac{\Delta t}{2} \right) \Delta \vec{u}_i^n = & \left[\vec{H}_i^n - \left(\frac{\delta \vec{F}^+}{\delta x} \right)_{i-\frac{1}{2}}^n - \left(\frac{\delta \vec{F}^-}{\delta x} \right)_{i+\frac{1}{2}}^n \right] + \frac{\Delta t}{2\delta x} \left\{ \left[\Psi^- \left(\delta \vec{F} - \Delta t \mathbf{J} \vec{G} \right)^- \right]_{i+\frac{3}{2}}^n \right. \\ & \left. - \left[\Psi^- \left(\delta \vec{F} - \Delta t \mathbf{J} \vec{G} \right)^- \right]_{i+\frac{1}{2}}^n - \left[\Psi^+ \left(\delta \vec{F} + \Delta t \mathbf{J} \vec{G} \right)^+ \right]_{i-\frac{1}{2}}^n + \left[\Psi^+ \left(\delta \vec{F} + \Delta t \mathbf{J} \vec{G} \right)^+ \right]_{i-\frac{3}{2}}^n \right\} \quad (5.3) \end{aligned}$$

The flux limiter functions $\Psi(r)$ have to obey a series of conditions in order to preserve the second order. They are

$\Psi(r) = 1, \forall r \Rightarrow$ Second order in space and time upwind scheme

$\Psi(r) = r, \forall r \Rightarrow$ Lax-Wendroff scheme (5.4)

Infinite schemes can be created again that would be intermediate between the fully upwind or fully central second order approximations. The zone of second order in space and time schemes is shaded in Fig. 5.1.

For the second order in space and time scheme the TVD second order region is the same that of Fig. 4.2, so that all the limiting functions defined in the chapter 4 are also suitable for this scheme.

The simplest way to extend the second order in space and time TVD scheme to systems with source terms is to start with the shorthand notation

$$\vec{R}^\pm = (\delta \vec{F} \pm \Delta t \mathbf{J} \vec{G})^\pm$$

with $\vec{R} = (R^1, \dots, R^k)$, and to define the flux limiting functions Ψ^+ and Ψ^- as matrices of the form

$$\Psi_{i+\frac{1}{2}}^\pm = \begin{pmatrix} \Psi \left(\frac{(R^1)^\pm_{i+\frac{1}{2} \pm 1}}{(R^1)^\pm_{i+\frac{1}{2}}} \right) & & \\ & \ddots & \\ & & \Psi \left(\frac{(R^k)^\pm_{i+\frac{1}{2} \pm 1}}{(R^k)^\pm_{i+\frac{1}{2}}} \right) \end{pmatrix}$$

The desired second order properties are then achieved and, although the TVD condition is only strictly ensured for the characteristic variables, it produces minimum oscillations in practice when

dealing with conservative variables. The stability condition is $CFL \leq 1$ loosing by limiting the second order upwind scheme the property of remaining stable up to $CFL = 2$.

The scheme (5.3) defined is conservative since it admits a wave decomposition like

$$\begin{aligned}\vec{F}_i^T &= \vec{F}_i^n \\ \delta \vec{F}_{i+\frac{1}{2}}^L &= \left(\delta \vec{F}^+ \right)_{i+\frac{1}{2}}^n - \frac{1}{2} (\Psi^+ \vec{R}^+)_{i-\frac{1}{2}}^n + \frac{1}{2} (\Psi^- \vec{R}^-)_{i+\frac{3}{2}}^n \\ \delta \vec{F}_{i+\frac{1}{2}}^R &= \left(\delta \vec{F}^- \right)_{i+\frac{1}{2}}^n + \frac{1}{2} (\Psi^+ \vec{R}^+)_{i-\frac{1}{2}}^n - \frac{1}{2} (\Psi^- \vec{R}^-)_{i+\frac{3}{2}}^n\end{aligned}$$

and a cell numerical flux

$$\vec{F}_{i+\frac{1}{2}}^* = \frac{1}{2} \left\{ \vec{F}_i^n + \vec{F}_{i+1}^n - \left[\mathbf{P} \text{sign}(\Lambda) \mathbf{P}^{-1} \delta \vec{F} \right]_{i+\frac{1}{2}}^n + (\Psi^+ \vec{R}^+)_{i-\frac{1}{2}}^n - (\Psi^- \vec{R}^-)_{i+\frac{3}{2}}^n \right\}$$

Note which the vector \vec{R} includes source terms, therefore a coupling of flux and source terms appears in the second order in time terms.

In order to complete the construction of the scheme, an upwind treatment of the source terms will be applied now. For that purpose, and using again the notation based on variable $\vec{G}_{i+\frac{1}{2}}$, the upwind TVD second order in space and time scheme can be expressed as a sum of the first order upwind plus limited corrections to second order as

$$\begin{aligned} \left(1 - \mathbf{K}_i^n \frac{\Delta t}{2} \right) \Delta \vec{u}_i^n &= \Delta t \left[\left(\vec{G}^+ \right)_{i-\frac{1}{2}}^n + \left(\vec{G}^- \right)_{i+\frac{1}{2}}^n \right] + \frac{\Delta t}{2} \left\{ \left[\Psi^+ \left(1 - \frac{\Delta t}{\delta x} \mathbf{J}^+ \right) \vec{G}^+ \right]_{i-\frac{1}{2}}^n \right. \\ &\quad \left. - \left[\Psi^+ \left(1 - \frac{\Delta t}{\delta x} \mathbf{J}^+ \right) \vec{G}^+ \right]_{i-\frac{3}{2}}^n + \left[\Psi^- \left(1 + \frac{\Delta t}{\delta x} \mathbf{J}^- \right) \vec{G}^- \right]_{i+\frac{1}{2}}^n - \left[\Psi^- \left(1 + \frac{\Delta t}{\delta x} \mathbf{J}^- \right) \vec{G}^- \right]_{i+\frac{3}{2}}^n \right\} \quad (5.5) \end{aligned}$$

And, Lax-Wendroff scheme, on the other hand, with upwind source term can be expressed

$$\begin{aligned} \left(1 - \mathbf{K}_i^n \frac{\Delta t}{2} \right) \Delta \vec{u}_i^n &= \Delta t \left[\left(\vec{G}^+ \right)_{i-\frac{1}{2}}^n + \left(\vec{G}^- \right)_{i+\frac{1}{2}}^n \right] + \frac{\Delta t}{2} \left\{ \left[\left(1 - \frac{\Delta t}{\delta x} \mathbf{J}^+ \right) \vec{G}^+ \right]_{i+\frac{1}{2}}^n \right. \\ &\quad \left. - \left[\left(1 - \frac{\Delta t}{\delta x} \mathbf{J}^+ \right) \vec{G}^+ \right]_{i-\frac{1}{2}}^n + \left[\left(1 + \frac{\Delta t}{\delta x} \mathbf{J}^- \right) \vec{G}^- \right]_{i-\frac{1}{2}}^n - \left[\left(1 + \frac{\Delta t}{\delta x} \mathbf{J}^- \right) \vec{G}^- \right]_{i+\frac{1}{2}}^n \right\} \quad (5.6) \end{aligned}$$

Defining now

$$\vec{S}^+ = \left(1 - \frac{\Delta t}{\delta x} \mathbf{J}^+ \right) \vec{G}^+, \quad \vec{S}^- = \left(1 + \frac{\Delta t}{\delta x} \mathbf{J}^- \right) \vec{G}^-$$

In order to obtain with $\Psi(r) = r$ the Lax-Wendroff scheme (5.6), for preserving the second order of approximation, the diagonal limiting matrices have to be

$$\Psi_{i+\frac{1}{2}}^{\pm} = \begin{pmatrix} \Psi \left(\frac{(S^1)^{\pm}_{i+\frac{1}{2} \pm 1}}{(S^1)^{\pm}_{i+\frac{1}{2}}} \right) & & \\ & \ddots & \\ & & \Psi \left(\frac{(S^k)^{\pm}_{i+\frac{1}{2} \pm 1}}{(S^k)^{\pm}_{i+\frac{1}{2}}} \right) \end{pmatrix}$$

6 Numerical results

A set of tests has been selected to illustrate the performance of some of the techniques described in the paper. The examples applies to the system of shallow water equations. In all of them, the primitive version of the conservative schemes (2.8) has been applied. When source terms are present, the semi-implicit upwind treatment has been implemented with $\theta = 0.5$. The numerical treatment at the boundaries follows the lines described into the appendix B and the entropy correction used in all the schemes is described into the appendix A.

6.1 Dam-break flow

This classical test case is considered a benchmark for comparison of the performance of numerical schemes specially designed for discontinuous transient flow. It allows to go a step further since the dam-break problem is defined by the non-linear system of homogeneous shallow water equations. Starting from initial conditions given by still water and two different water levels separated by a dam, the theory of characteristics supplies an exact evolution solution [9] that can be used as a reference. In the example presented, two ratios of initial water depths $\frac{h_L}{h_R} = 10$ and $\frac{h_L}{h_R} = 100$ are used. The solution is displayed in Figs. 6.1-6.6 for $t = 20s$. The CFL number used is 90% of the maximum for stability. In spatial second order TVD scheme this is 0.45 with Superbee flux limiter and 0.6 with Minmod flux limiter, being 0.9 in the others schemes. An interval of $\Delta x = 1m$ is used in the mesh. The entropy correction described into appendix A is applied in all cases and it is remarkable that it makes a good results, being the typical "dog-leg" effect negligible. It is also remarkable that the Lax-Wendroff scheme with entropy correction, although appreciable numerical oscillations are shown, is able to solve strong shocks without to be necessary a TVD correction.

First order upwind scheme provides a reasonably good results with a slight numerical diffusion. The second order in space TVD scheme tends to produce antidiffusive solutions, being this very evident with the Superbee flux limiter. Nevertheless with the Minmod flux limiter this is less antidiffusive providing a slight improvement with regard to the first order scheme. Second order in space and time improves the numerical solution being the more accurate scheme.

6.2 Steady flow in channels

When the shallow water equations are used to model hydraulic problems involving bed slope changes and bed friction, the system is no longer homogeneous and the source terms have to be taken into account. On the other hand, this renders more difficult and often impossible to find reference exact solutions. McDonald [7] proposed a set of test cases based on steady flow in channels of varying bed slope and/or breadth by calculating the analytical slope and breadth functions compatible with constant discharge conditions given an analytical water depth function. Among them, we have chosen an example consisting of a 650m long trapezoidal channel with a bed variation given by a slope function of x and a roughness coefficient $n = 0.03$. The constant discharge $Q = 20m^3/s$ is imposed upstream and supercritical flow is enforced downstream. There are some points of transcritical flow. The same CFL number to the dam-break problem is used in this case. The interval of the mesh is $\Delta x = 6.5m$. The most important detail to note on the results is the perfect coincidence of all the upwind based schemes with upwind treatment of the source terms and the limitation proposed in the TVD schemes versus the results given by the Lax-Wendroff scheme (Fig. 6.16). This is an expectable behaviour since it is a steady flow problem in which the upwind treatment of the source terms produces second order spatial accuracy even in the first order upwind scheme as mentioned before. In Fig. 6.12 the "traditional" form of the TVD scheme (with pointwise and without limiting treatment of the source terms) is performed being evident the better results with an upwind and limited treatment of the source terms. In Fig. 6.11 two different forms of limiting the source terms are compared. A slight loss of accuracy excluding the source terms in the flux limiter is patent in the discontinuities because it is the place where the flux limiter plays a determining role.

6.3 Unsteady flow in rivers

In order to show the application to a practical case, an example of unsteady flow in a river is presented now. It is a 9000m long reach of the upstream part of river Neila in Spain. Being a mountain river, it is characterized by strong irregularities in the cross section, by a rather steep part in the first kilometres and by a low base discharge ($1m^3/s$) which, altogether, produce a high velocity basic flow, transcritical in some parts. In order to check the conservation properties of the schemes applied, and the absence of oscillations in the TVD schemes, a sudden increase in discharge to $40m^3/s$ and a critical depth is imposed at the upstream end. This step hydrograph propagates into the river. The same CFL number as the steady flow cases and an interval of $\Delta x = 22.5m$ in the mesh are used. Figs. 6.13-6.16 show that it does so with almost a perfectly constant value at times $t = 500s$, $t = 1000s$ and $t = 1500s$. In Fig. 6.17 it is seen the oscillations produced by non-limiting the source terms. Fig. 6.18 shows the detail of the front wave where the advantages of using higher order approaches are remarkable, this is not so clear when reproducing steady states. In Fig. 6.19 the strong gradient in the bed slope of river Neila can be seen. Figs. 6.20-6.24 show some other parameters with the second order in space and time TVD scheme with Superbee limiter (the most accurate scheme) and the strong irregularities of the river are evident.

7 Summary and conclusions

In this work a study of different one step explicit schemes is presented. The one step schemes, although slightly more complex than two step methods (eg McCormack [5,6]) of second order of approximation, are faster for resolving shallow water equations in rivers and irregular channels because the time elapsed calculating the sectional parameters of channel is much greater than the time elapsed resolving the numerical schemes. Predictor-corrector schemes calculate twice the parameters every time step requiring double CPU time versus one step schemes. Moreover, one step methods admit a semi-implicit treatment of source terms, necessary for stabilising the simulations when these become dominant (rivers and irregular channels), whereas two step schemes loose the second order of approximation in time property when a semi-implicit treatment of source terms is applied.

The applications presented are based on conservative schemes in the primitive form. This form takes advantage of the simplicity of the primitive form of the equations and produces faster and

simpler schemes without losing the accuracy of conservative schemes.

The first order upwind scheme with upwind and semi-implicit treatment of the source terms, at this moment, is one of the best schemes simulating shallow water equations because, although this is a first order scheme, it is robust, reasonably simple, fast and it produces second order solutions in steady and quasi-steady problems. Therefore this is the preferential explicit scheme in steady flows. On the other hand, it is very accurate solving strong shocks which many other schemes cannot simulate. The conservation errors are very small and absence of oscillations is seen. Nevertheless, the results are less accurate in unsteady flows.

The new entropy correction proposed into the appendix A for the upwind schemes improves the solutions when the accuracy of this correction is crucial. The "*dog-leg*" effect is negligible as it is shown in dam-break tests.

If maximal accuracy in unsteady flows is required, one of the second order TVD schemes is recommended. In these schemes, an upwind and a limiting treatment of source terms produces better results, being insignificant the numerical oscillations produced in the propagated shocks, than a pointwise treatment of these terms, which a small oscillations and a bad conservation are seen. Furthermore, a slight improvement is achieved incorporating the source terms in the flux limiter.

The TVD second order in space scheme is equally robust, more accurate, but slightly more complex than the first order upwind scheme and is the simplest high order TVD scheme. It is recommended to use the Minmod flux limiter with this scheme because this achieves the best results and is less restrictive in the CFL condition ($CFL < \frac{2}{3}$ versus $CFL < \frac{1}{2}$). However, this limitation in the CFL number is compensated by an increment of 50% in the CPU time compared with all other methods used in this work.

Though a little more complex, the TVD second order in space and time scheme is faster (the CPU time is similar to the first order schemes) than the TVD second order in space scheme. With this scheme the Superbee and Van-Leer flux limiters produce more accurate solutions representing a small improvement with regard to the TVD second order in space scheme. Therefore, it is the

best of the analysed schemes presenting the best performance in unsteady flows. Nevertheless, this is a cost of a greater complexity and the results on shallow water equations are better than the first order upwind scheme for highly unsteady flows but identical for steady flows.

8 * References

- [1] Alcrudo, F. (1992). " *Esquemas de alta resolucion de variacion total decreciente para el estudio de flujos discontinuos de superficie libre* ". PhD thesis, Universidad de Zaragoza.
- [2] Casier, F., Deconinck, H., and Hirsch, C. (1984). " *A class of bidiagonal schemes for solving the Euler equations* ". **AIAA J.**, 22(11), 1556-1563.
- [3] Garcia Navarro P., and Hubbard M. E. (1999). " *Flux difference splitting and the balancing of source terms and flux gradients* ". Internal Report 3/99, Dep. of Mathematics, Univ. of Reading.
- [4] Garcia Navarro P., and Vazquez Cendon M. E. (1997). " *Some considerations and improvements on the performance of Roe's scheme for 1D irregular geometries* ". Internal Report 23, Submitted to Computers and Fluids. Dep. de Matematica Aplicada, Univ. de Santiago do Compostela .
- [5] Hirsch C. (1990). " *Computational methods for inviscid and viscous flows: Numerical computation of internal and external flows* ". Vol. 2, John Wiley & Sons, New York.
- [6] McCormack R. W. (1971). " *Numerical solutions of the interaction of a shock wave with a laminar boundary layer* ". **Lectures notes in Physics**, 81.
- [7] McDonald, I. (1996). " *Analysis and computation of steady open channel flow* ". PhD thesis, University of Reading.
- [8] Roe, P. L. (1981). " *Approximate Riemann solvers, parameter vectors, and difference schemes* ". **Journal of Computational Physics**, 43(2), 357-372.
- [9] Stoker, J. J. (1957). " *Water waves* ". Interscience Pub. New York.
- [10] Sweby, P. K. (1984). " *High resolution schemes using flux limiters for hyperbolic conservation laws* ". **SIAM J. Numer. Anal.**, 21 995-1011.

[11] Toro E. (1997). " *Riemann solvers and numerical methods for fluid dynamics: a practical introduction* ". Springer, Berlin.

Appendix A: Solution to the entropy problem with artificial viscosity

A.1 Scalar case

For the resolution of scalar propagation equation like

$$\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0, \quad a = a(x, u)$$

If there is a transition from subcritical ($a < 0$) to supercritical ($a > 0$) flow between two nodal points in the grid, the way in which the situation will evolve in one Δt can be studied. This is illustrated in Fig A.1

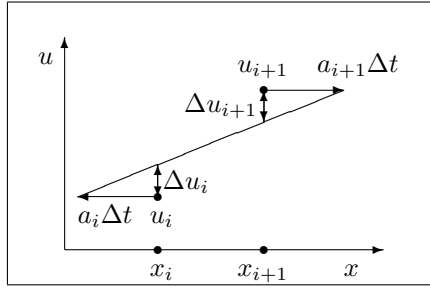


Figure A.1: Evolution of a propagation wave in a cell of the mesh with transcritical flux at an infinitesimal time Δt .

The value of the increments in the variable can be calculated by linear interpolation

$$\Delta u_i = -\frac{a_i \Delta t}{\delta x + (a_{i+1} - a_i) \Delta t} \delta u_{i+\frac{1}{2}} = -\frac{\sigma_i \delta u_{i+\frac{1}{2}}}{1 + \delta \sigma_{i+\frac{1}{2}}}$$

$$\Delta u_{i+1} = -\frac{a_{i+1} \Delta t}{\delta x + (a_{i+1} - a_i) \Delta t} \delta u_{i+\frac{1}{2}} = -\frac{\sigma_{i+1} \delta u_{i+\frac{1}{2}}}{1 + \delta \sigma_{i+\frac{1}{2}}}$$

with $\sigma = a \frac{\Delta t}{\delta x}$. These will be increments in first order. Then, the viscosity necessary in a numerical scheme to solve properly this kind of transitions can be found. Concentrating on the grid cell where the transition takes place, the conservative numerical schemes in general will follow a wave decomposition like

$$\delta F_{i+\frac{1}{2}}^L = a_{i+\frac{1}{2}}^L \delta u_{i+\frac{1}{2}}, \quad \delta F_{i+\frac{1}{2}}^R = a_{i+\frac{1}{2}}^R \delta u_{i+\frac{1}{2}}$$

$$\Delta u_{i+1} = -\sigma_{i+\frac{1}{2}}^L \delta u_{i+\frac{1}{2}}, \quad \Delta u_i = -\sigma_{i+\frac{1}{2}}^R \delta u_{i+\frac{1}{2}}$$

Adding artificial viscosity to fix the entropy problem

$$\Delta u_{i+1} = -(\sigma_{i+\frac{1}{2}}^L + \nu_{i+\frac{1}{2}})\delta u_{i+\frac{1}{2}}, \quad \Delta u_i = -(\sigma_{i+\frac{1}{2}}^R - \nu_{i+\frac{1}{2}})\delta u_{i+\frac{1}{2}}$$

Identifying these increments with those obtained by interpolation two possible values are found. To avoid problems, the maximum is chosen

$$\nu_{i+\frac{1}{2}} = \max \left(\sigma_{i+\frac{1}{2}}^R - \frac{\sigma_i}{1 + \delta\sigma_{i+\frac{1}{2}}}, \frac{\sigma_{i+1}}{1 + \delta\sigma_{i+\frac{1}{2}}} - \sigma_{i+\frac{1}{2}}^L \right)$$

This viscosity, having been obtained from a linear interpolation, gives a correction in first order, then it is valid for first order schemes. Furthermore, considering that there will be a finite number of transitions like this, it can be acceptable also for second order schemes. As the only first order scheme considered with entropy problems is the first order upwind, and the second order schemes can be expressed in terms of that first order scheme plus second order corrections, the previous solution is written as

$$\nu_{i+\frac{1}{2}} = \max \left(\sigma_{i+\frac{1}{2}}^- - \frac{\sigma_i}{1 + \delta\sigma_{i+\frac{1}{2}}}, \frac{\sigma_{i+1}}{1 + \delta\sigma_{i+\frac{1}{2}}} - \sigma_{i+\frac{1}{2}}^+ \right) \quad (\text{A.1})$$

Another forms of artificial viscosity to fix the entropy problem are described in [5].

A.2 Systems of equations

When dealing with homogeneous systems of equations, in order to study separately the behaviour of every wave, the system is first formulated in characteristic variables

$$\frac{\partial w^k}{\partial t} + a^k \frac{\partial w^k}{\partial x} = 0$$

The entropy problem can be fixed by analogy to the scalar case acting over every k component in the decoupled system. The artificial viscosity is then defined as a diagonal matrix. For the particular case of having a transition sub-super between nodes i and $i + 1$ in the k component

$$(\mathbf{V}^{kk})_{i+\frac{1}{2}} = \max \left((\sigma^k)_{i+\frac{1}{2}}^- - \frac{(\sigma^k)_i}{1 + \delta(\sigma^k)_{i+\frac{1}{2}}}, \frac{(\sigma^k)_{i+1}}{1 + \delta(\sigma^k)_{i+\frac{1}{2}}} - (\sigma^k)_{i+\frac{1}{2}}^+ \right) \quad (\text{A.2})$$

with $\sigma^k = \frac{\Delta t}{\delta x} a^k$. The characteristic system then is transformed to

$$\frac{\partial \vec{w}}{\partial t} + \mathbf{\Lambda} \frac{\partial \vec{w}}{\partial x} = \mathbf{V} \frac{\partial^2 \vec{w}}{\partial x^2}$$

Returning to the physical variables by means of \mathbf{P} y \mathbf{P}^{-1} , the matrices that diagonalize the Jacobian \mathbf{J} :

$$\mathbf{P} \left(\frac{\partial \vec{w}}{\partial t} + \mathbf{A} \frac{\partial \vec{w}}{\partial x} = \mathbf{V} \frac{\partial^2 \vec{w}}{\partial x^2} \right) \Rightarrow \frac{\partial \vec{u}}{\partial t} + \mathbf{J} \frac{\partial \vec{u}}{\partial x} = \mathbf{PVP}^{-1} \frac{\partial^2 \vec{u}}{\partial x^2}$$

and the artificial viscosity is defined as

$$\nu_{i+\frac{1}{2}} = (\mathbf{PVP}^{-1})_{i+\frac{1}{2}} \quad (\text{A.3})$$

with \mathbf{V} defined in (A.2).

As a last remark, the entropy problem does only affect a transition in fluxes and is not modified by the presence of source terms. Hence, in presence of source terms, (A.3) will still be used.

Appendix B: Boundary conditions

For the treatment at the boundaries, the distinction between numerical boundary conditions and external (imposed) boundary conditions has been exploited. The use of the characteristic variables to obtain the correct region of dependence of a point is a suitable method to produce numerical boundary conditions [2,5]. In this work, however, a similar way to generate numerical boundary conditions for any conservative scheme keeping the degree of accuracy of the scheme has been developed. All the conservative schemes admit a decomposition in sum of contributions from left and right like

$$\Delta \vec{u}_i^n = \Delta t \left(\vec{G}_{i+\frac{1}{2}}^R + \vec{G}_{i-\frac{1}{2}}^L \right) \quad (\text{B.1})$$

A decentralised treatment of source terms is supposed but a pointwise treatment is also possible.

At the boundaries, the following information can be used

$$\text{Upstream: } \Delta \vec{u}_i^n = \Delta t \vec{G}_{i+\frac{1}{2}}^R, \quad \text{Downstream: } \Delta \vec{u}_i^n = \Delta t \vec{G}_{i-\frac{1}{2}}^L$$

Therefore, the numerical boundary conditions can be worked out from the scheme itself whenever the domain of dependence of the boundary points is inside the calculation grid. Otherwise, physical or imposed external information must be used. Let us call $\Delta \vec{u}_i^N$ the numerical increments of the variable obtained using (B.1), and $\Delta \vec{u}_i^F$ the physical increments, that is, the final values that we want to determine. Then, if w^k is the characteristic variable associated to a propagation velocity a^k

that requires a numerical boundary condition, the following will be used

$$(\Delta w^k)_i^F = (\Delta w^k)_i^N \quad (\text{B.2})$$

The procedure to implement this technique in the shallow water equations requires first to know the kind of flow at every boundary. Then, (B.1) is used to get the numerical increments $\Delta \vec{u}_i^N$. Once calculated, they are used in (B.2) to obtain two possible numerical boundary conditions, one associated to the velocity $v + c$ and the other associated to the velocity $v - c$. The characteristic variables in this problem can be defined from

$$\Delta \vec{w} = \frac{1}{2c} \begin{pmatrix} (c - v)\Delta A + \Delta Q \\ (c + v)\Delta A - \Delta Q \end{pmatrix}$$

And there are four possibilities:

1. Subcritical inlet: One physical and one numerical boundary condition

$$(c + v)_i^n \Delta A_i^F - \Delta Q_i^F = (c + v)_i^n \Delta A_i^N - \Delta Q_i^N$$

2. Supercritical inlet: Two physical boundary conditions

$$\Delta A_i^F = datum, \quad \Delta Q_i^F = datum$$

3. Subcritical outlet: One physical and one numerical boundary condition

$$(c - v)_i^n \Delta A_i^F + \Delta Q_i^F = (c - v)_i^n \Delta A_i^N + \Delta Q_i^N$$

4. Supercritical outlet: Two numerical boundary conditions

$$\Delta A_i^F = \Delta A_i^N, \quad \Delta Q_i^F = \Delta Q_i^N$$